

I claim:

1. A computer-assisted method for identifying duplicate and near-duplicate documents in a large collection of documents, comprising the steps of:

initially, selecting distinctive features contained in the collection of documents,

5 then, for each document, identifying the distinctive features contained in the document, and

then, for each pair of documents having at least one distinctive feature in common, comparing the distinctive features of the documents to determine whether the documents are duplicate or near-duplicate documents.

2. The computer-assisted method according to claim 1, wherein the method is applied to removing duplicates in document collections.

3. The computer-assisted method according to claim 1, wherein the method is applied to detecting plagiarism.

4. The computer-assisted method according to claim 1, wherein the method is applied to detecting copyright infringement.

5. The computer-assisted method according to claim 1, wherein the method is applied to determine the authorship of a document.

6. The computer-assisted method according to claim 1, wherein the method is applied to clustering successive versions of a document from among a collection of documents.

7. The computer-assisted method according to claim 1, wherein the method is applied to seeding a text classification or text clustering algorithm with sets of duplicate or near-duplicate documents.

8. The computer-assisted method according to claim 1, wherein the method is applied to matching an e-mail message with responses to the e-mail message.

9. The computer-assisted method according to claim 1, wherein the method is applied to matching responses to an e-mail message with the e-mail message.

10. The computer-assisted method according to claim 1, wherein the method is applied to creating a document index for use with a query system to efficiently find documents in response to a query which contain a particular phrase or excerpt.

11. The computer-assisted method according to claim 10, wherein the document index can be utilized even if the particular phrase or excerpt was not recorded correctly in the document or in the query.

12. The computer-assisted method according to claim 1, wherein the distinctive features appear in a different order in each of the documents.

13. The computer-assisted method according to claim 1, wherein the distinctive features are distinctive text fragments from the documents in the document collection.

14. The computer-assisted method according to claim 13, wherein the method is applied to information retrieval methods.

15. The computer-assisted method according to claim 14, wherein the information retrieval method is a text classification method.

16. The computer-assisted method according to claim 14, wherein:
the information retrieval method assumes word independence, and
the distinctive text fragments are added to an index set.

17. The computer-assisted method according to claim 13, wherein the distinctive text fragments are sequences of at least two words that appear in a limited number of documents in the document collection.

18. The computer-assisted method according to claim 14, wherein if one distinctive text fragment is contained within another distinctive text fragment within the same document, only the longest distinctive text fragment is considered as a distinctive feature.

19. The computer-assisted method according to claim 17, wherein the sequences of at least two words are considered as appearing in a document when the document contains the sequence of at least two words at least a user-specified minimum number of times.

20. The computer-assisted method according to claim 17, wherein the sequences of at least two words are considered as appearing in a document when the document contains the sequence of at least two words at least a user-specified minimum frequency.

21. The computer-assisted method according to claim 17, wherein:
for each sequence of at least two words, a distinctiveness score is calculated, and
the highest scoring sequences that are found in at least two documents in
5 the document collection are considered distinctive text fragments.

22. The computer-assisted method according to claim 21, wherein the distinctiveness score is the reciprocal of the number of documents containing the phrase multiplied by a monotonic function of the number of words in the phrase.

23. The computer-assisted method according to claim 21, wherein the monotonic function is the number of words in the phrase.

24. The computer-assisted method according to claim 21, wherein the distinctiveness score is the percentage of documents not containing the phrase multiplied by a monotonic function of the number of words in the phrase.

25. The computer-assisted method according to claim 24, wherein the monotonic function is the number of words in the phrase.

26. The computer-assisted method according to claim 17, wherein the limited number is selected by a user.

27. The computer-assisted method according to claim 17, wherein the limited number is defined by a linear function of the number of documents in the document collection.

28. The computer-assisted method according to claim 17, wherein the distinctive text fragments include glue words.

29. The computer-assisted method according to claim 28, wherein the glue words do not appear at either extreme of the distinctive text fragments.

30. The computer-assisted method according to claim 1, further including the step of for each pair of documents having at least one distinctive feature in common, counting the number of distinctive features in common,

5 wherein determining whether the pair of documents is duplicates or near-duplicates includes the steps of:

for each pair of documents, calculating an overlap ratio by dividing the number of distinctive features in common by the smaller of the number of distinctive features per document, and

10 comparing the overlap ratio to a threshold and if the overlap ratio is greater than the threshold, then the pair of documents are duplicates or near-duplicates, otherwise the pair of documents is not duplicates or near-duplicates.

31. The computer-assisted method according to claim 30, further including the steps of:

building a document index that maps each document to its associated distinctive features, wherein if one distinctive feature is repeated within one document,
5 the index maps the document to the distinctive feature once, and

building a feature index that maps each distinctive feature to its associated document, wherein if one distinctive feature is repeated within one document, the index maps the distinctive feature to the document once,

wherein determining whether the pair of documents are duplicates or
10 near-duplicates further includes the steps of:

creating a list of unique distinctive features from the document index,

for each unique distinctive feature, creating a list of documents which contain the unique distinctive feature, and

15 for each document, creating a list of documents that have at least one feature in common with the document and the number of features in common with the document.

32. The computer-assisted method according to claim 31, wherein the distinctive features include distinctive phrases.

33. The computer-assisted method according to claim to 31, wherein the distinctive features appear in a different order in each of the documents.

34. The computer-assisted method according to claim 31, wherein the distinctive features include text spans.

35. The computer-assisted method according to claim 34, wherein the text spans include sentences.

36. The computer-assisted method according to claim 34, wherein the text spans include lines of text.

37. A computer-assisted method for identifying duplicate and near-duplicate text spans in a large collection of text spans, comprising the steps of:

initially, selecting distinctive features contained in the collection of text spans,

5 then, for each text span, identifying the distinctive features contained in the text span, and

then, for each pair of text spans having at least one distinctive feature in common, comparing the distinctive features of the text spans to determine whether the text spans are duplicate or near-duplicate text spans.

38. The computer-assisted method according to claim 37, wherein the text spans are sentences.

39. An apparatus to enable a method for identifying duplicate and near-duplicate documents in a large collection of documents, comprising:

a means for initially selecting distinctive features contained in the collection of documents;

5 a means for subsequently identifying the distinctive features contained in each document; and

a means for then comparing the distinctive features of each pair of documents having at least one distinctive feature in common to determine whether the documents are duplicate or near-duplicate documents.